

# Introduction to Big Data

Thomas Heinis

Imperial College  
London

# Topics

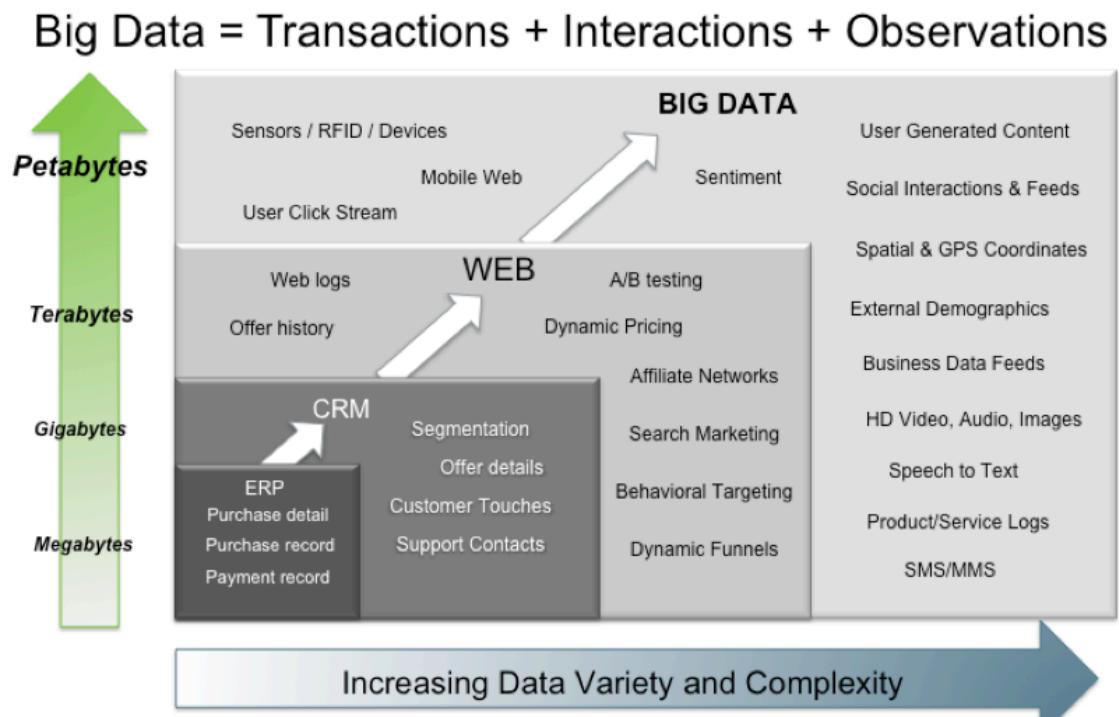
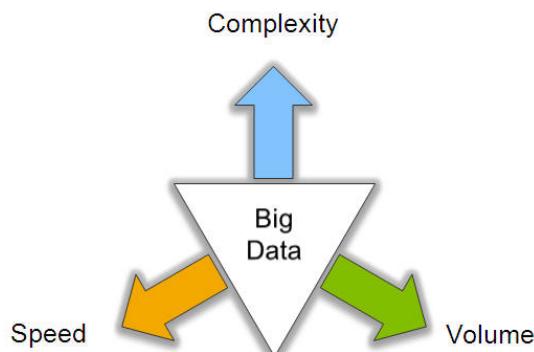
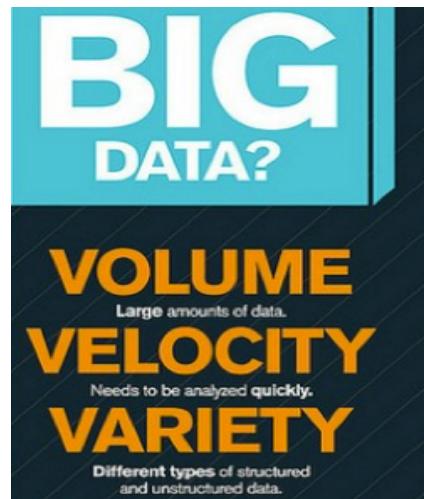
- Scope: Big Data & Analytics
- Topics:
  - Foundation of Scalable Data Management
  - Hadoop/Map-Reduce Programming and Data Processing & BigTable/Hbase/Cassandra
  - Graph Database and Graph Analytics

# What's Big Data?

**No single definition; from Wikipedia:**

- **Big data** is the term for a collection of data sets so large and complex that it becomes difficult to process using on-hand database management tools or traditional data processing applications.
- The challenges include capture, curation, storage, search, sharing, transfer, analysis, and visualization.
- The trend to larger data sets is due to the additional information derivable from analysis of a single large set of related data, as compared to separate smaller sets with the same total amount of data, allowing correlations to be found to “spot business trends, determine quality of research, prevent diseases, link legal citations, combat crime, and determine real-time roadway traffic conditions.”

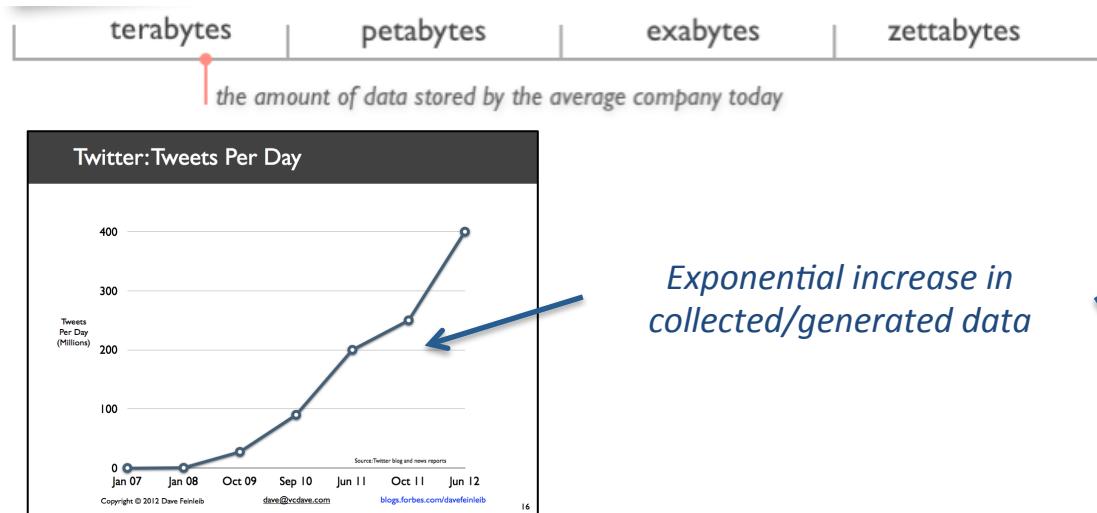
# Big Data: 3V's



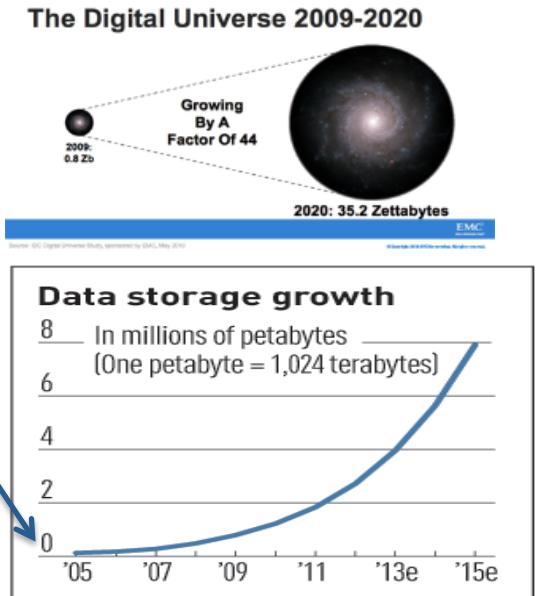
*Source:* Contents of above graphic created in partnership with Teradata, Inc.

# Volume (Scale)

- **Data Volume**
  - 44x increase from 2009 to 2020
  - From 0.8 zettabytes to 35zb
- Data volume is increasing exponentially



*Exponential increase in collected/generated data*



**12+ TBs**  
of tweet data  
every day



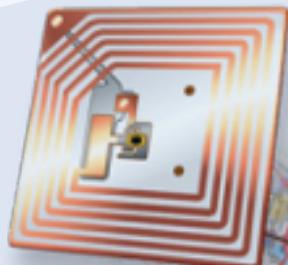
? TBs of  
data every day



**25+ TBs of**  
log data every  
day



**30 billion** RFID  
tags today  
(1.3B in 2005)



**76 million** smart meters  
in 2009...  
200M by 2014

http://

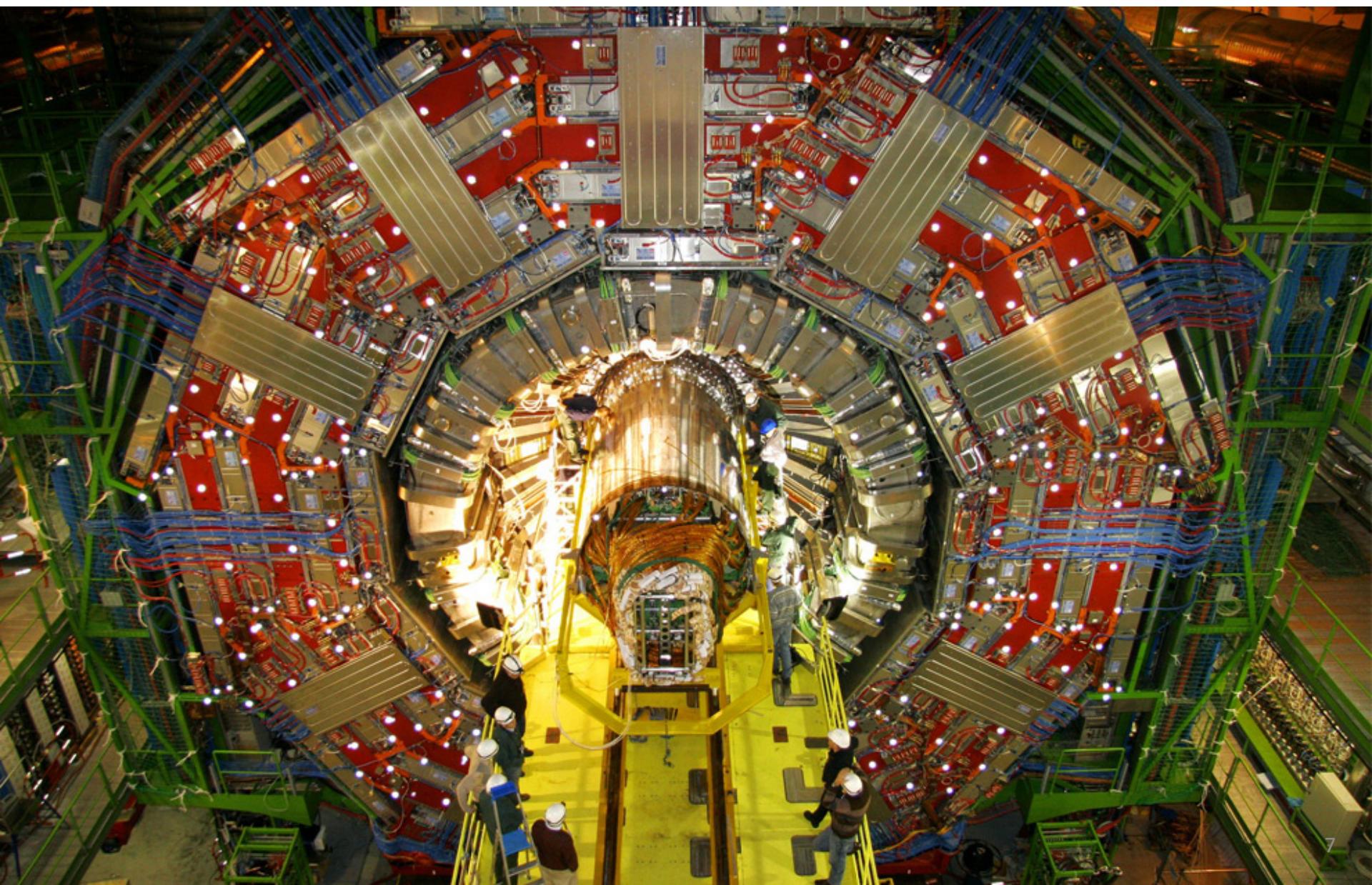
**4.6 billion**  
camera  
phones  
world wide

**100s of**  
**millions**  
**of GPS**  
**enabled**  
devices sold  
annually



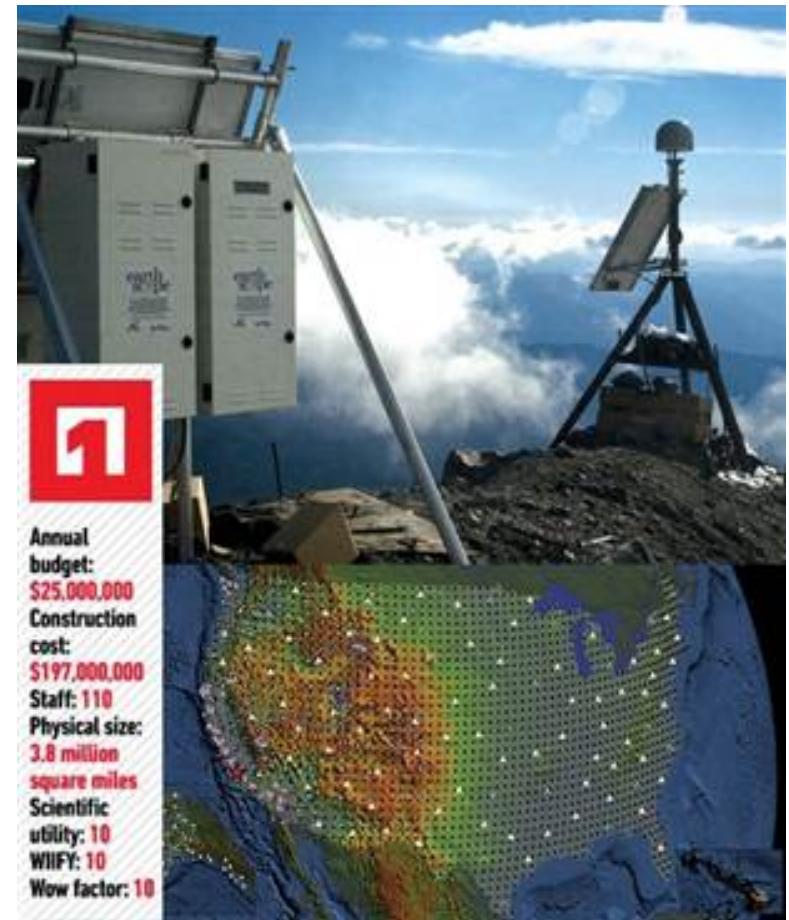
**2+**  
**billion**  
people on  
the Web  
by end  
2011

# CERN's Large Hydron Collider (LHC) generates 15 PB a year

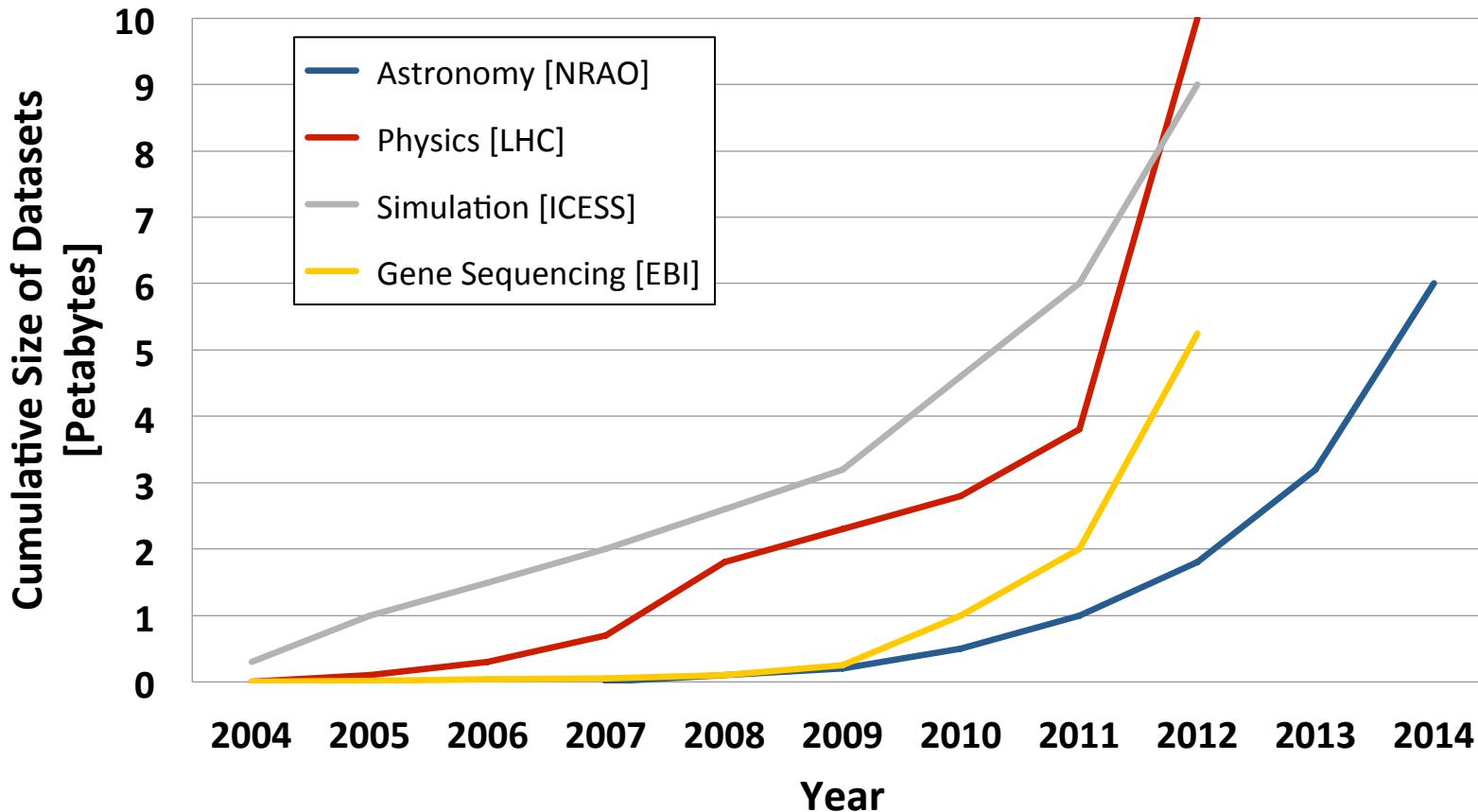


# The Earthscope

The Earthscope is the world's largest science project. Designed to track North America's geological evolution, this observatory records data over 3.8 million square miles, amassing 67 terabytes of data. It analyzes seismic slips in the San Andreas fault, sure, but also the plume of magma underneath Yellowstone and much, much more.



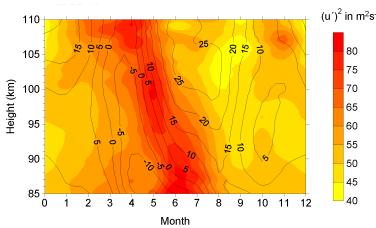
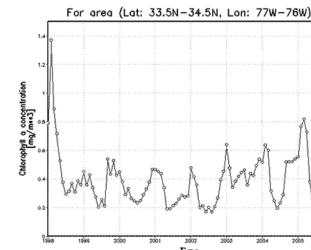
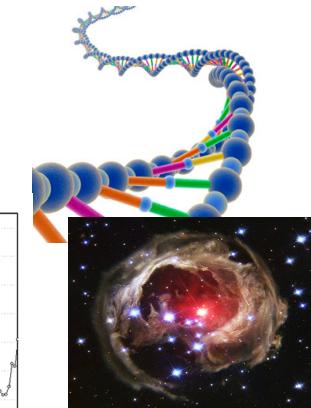
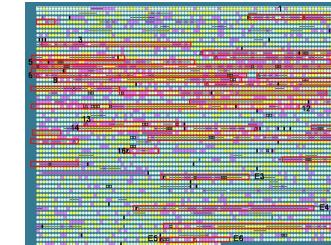
# Scientific Data Growth



Scientific Data Grows Exponentially!

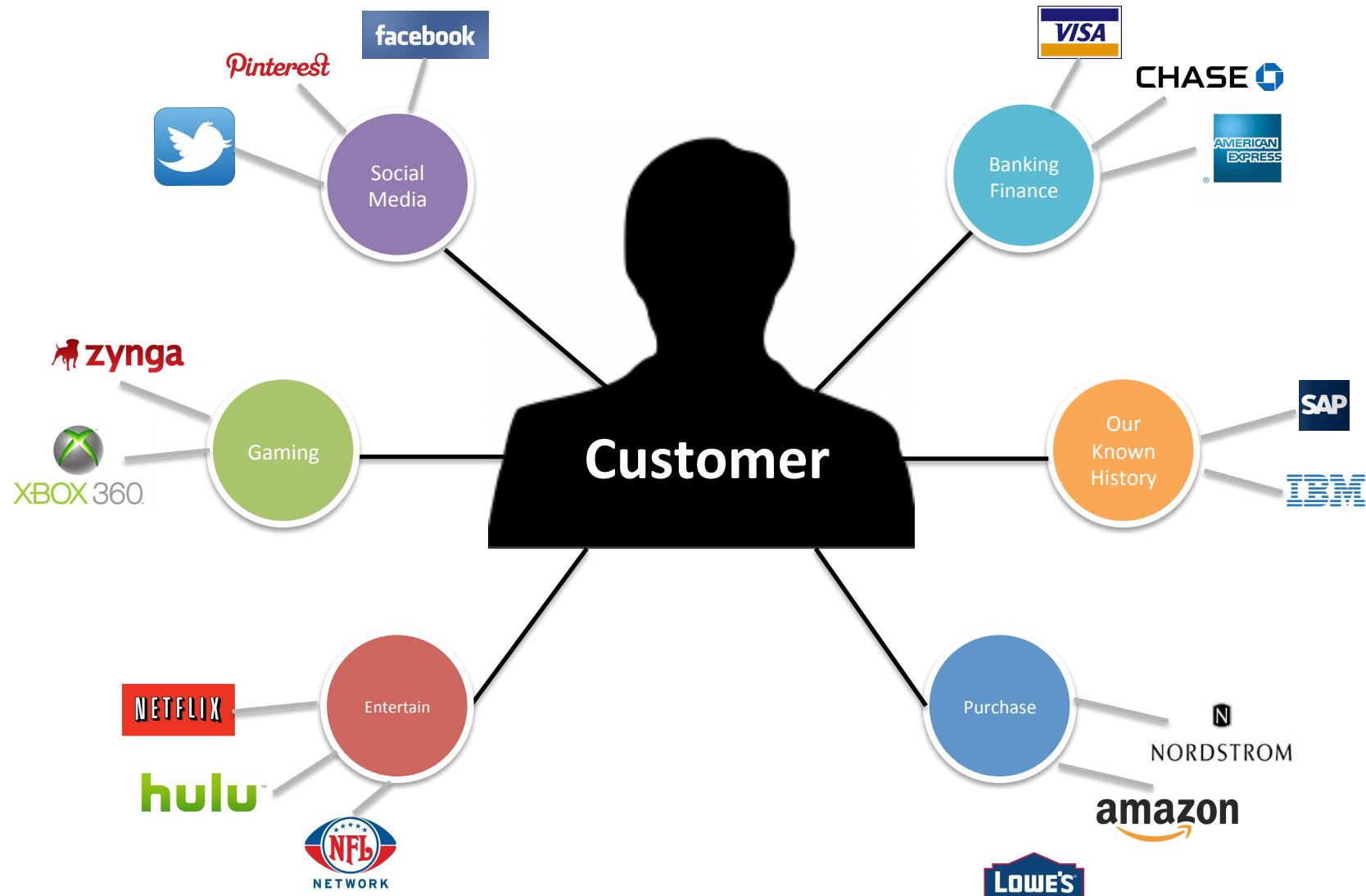
# Variety (Complexity)

- Relational Data (Tables/ Transaction/Legacy Data)
- Text Data (Web)
- Semi-structured Data (XML)
- Graph Data
  - Social Network, Semantic Web (RDF), ...
- Streaming Data
  - You can only scan the data once
- A single application can be generating/collecting many types of data
- Big Public Data (online, weather, finance, etc.)



To extract knowledge →  
All these types of data need to be linked together

# A Single View to the Customer



# Velocity (Speed)

- Data is being generated fast and need to be processed fast
- Online Data Analytics
- Late decisions → missing opportunities
- Examples
  - **E-Promotions:** Based on your current location, your purchase history, what you like → send promotions right now for store next to you
  - **Healthcare Monitoring:** sensors monitoring your activities and body → any abnormal measurements require immediate reaction



# Real-Time/Fast Data



**Social media and networks**  
(all of us are generating data)



**Scientific instruments**  
(collecting all sorts of data)



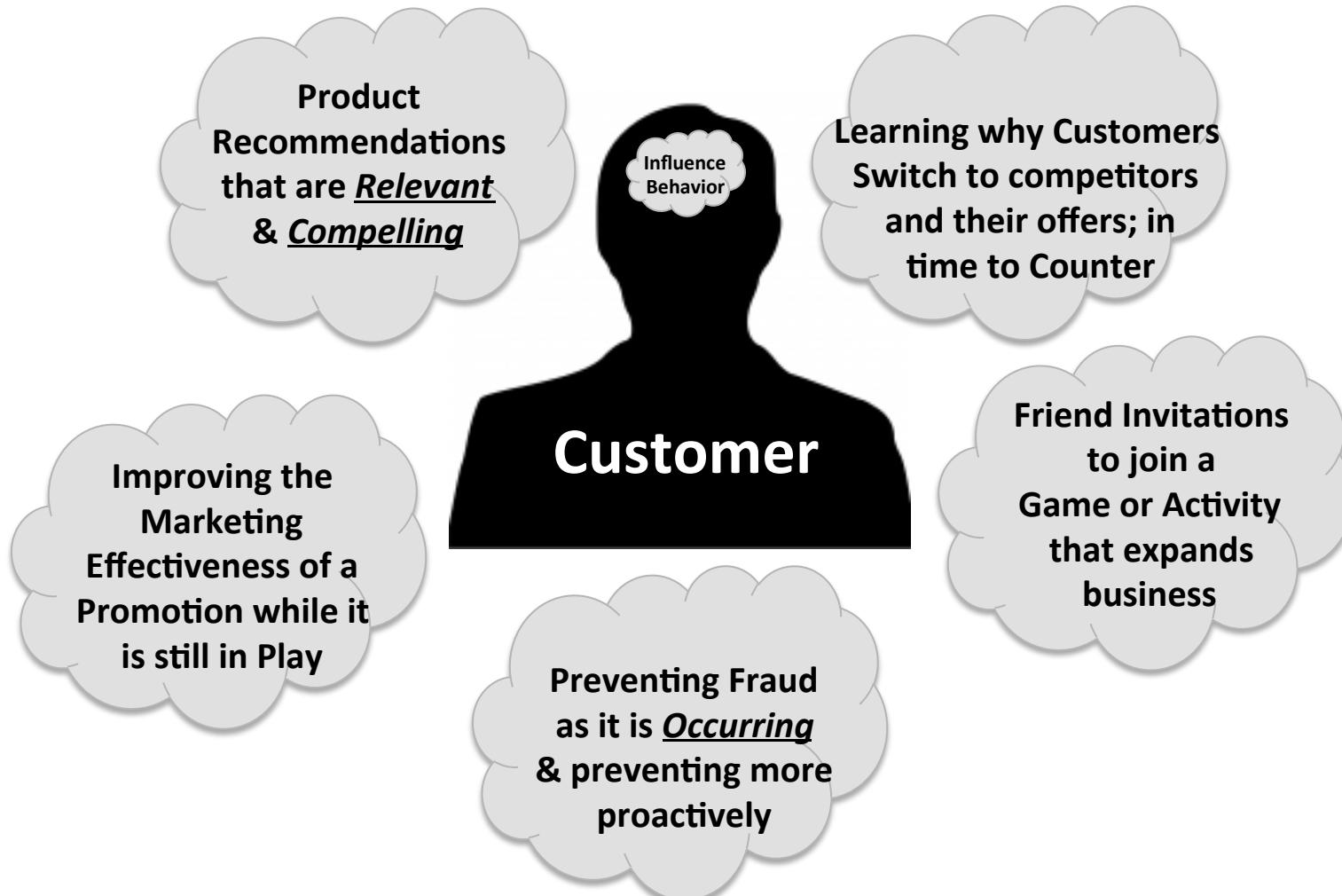
**Mobile devices**  
(tracking all objects all the time)



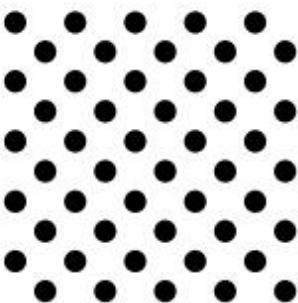
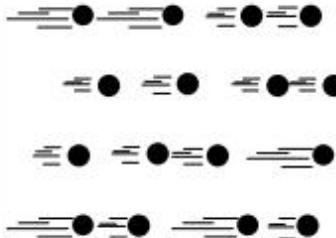
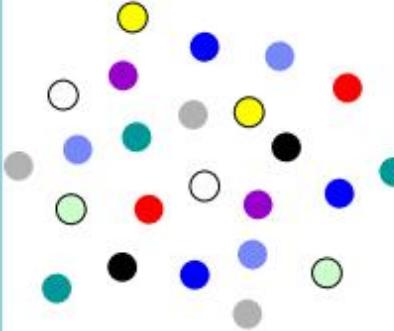
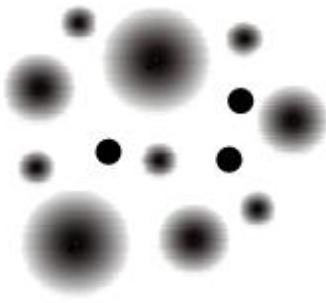
**Sensor technology and networks**  
(measuring all kinds of data)

- The progress and innovation is no longer hindered by the ability to collect data
- But, by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion

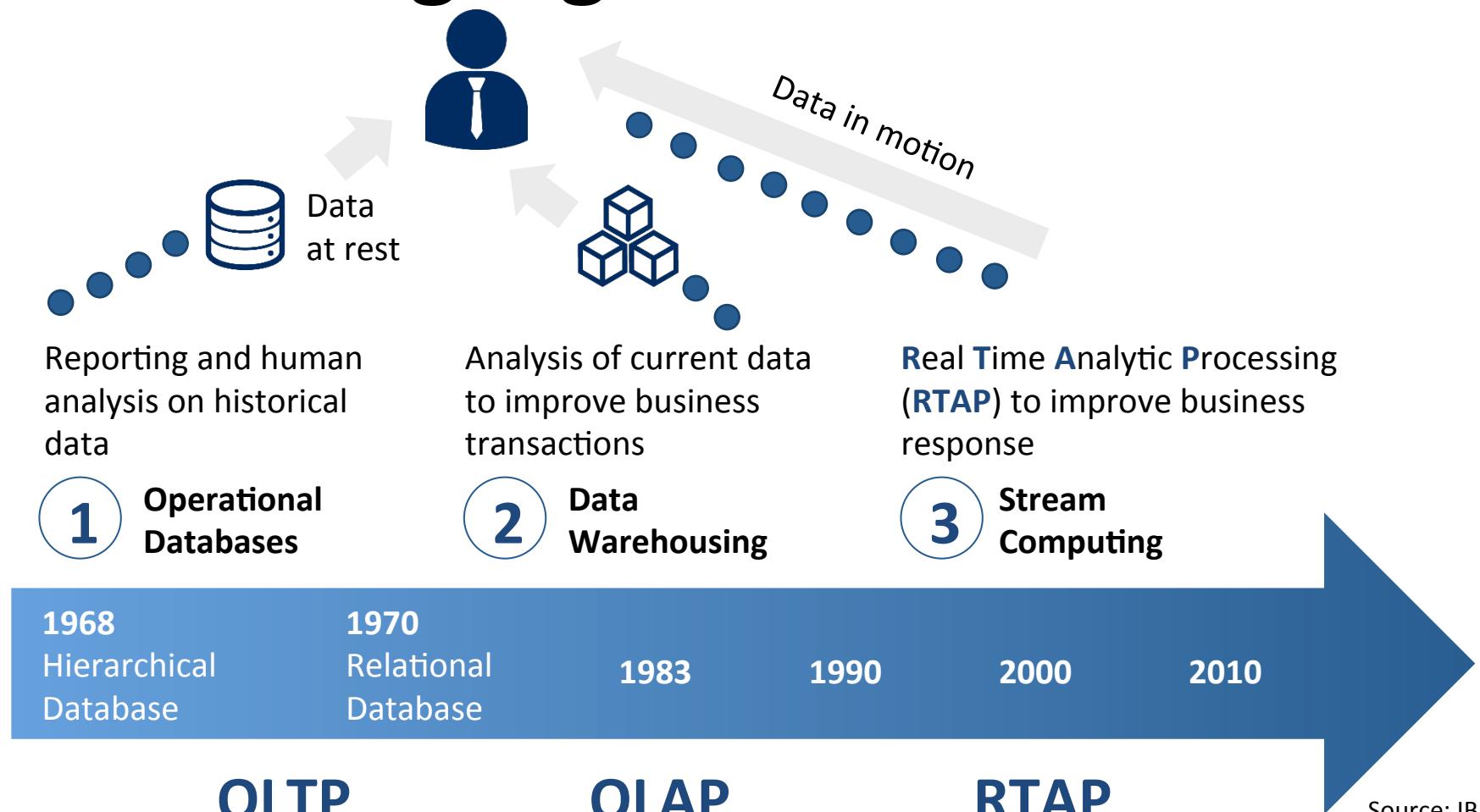
# Real-Time Analytics/Decision Requirement



# Some Make it 4V's

Volume	Velocity	Variety	Veracity*
			
<b>Data at Rest</b> Terabytes to exabytes of existing data to process	<b>Data in Motion</b> Streaming data, milliseconds to seconds to respond	<b>Data in Many Forms</b> Structured, unstructured, text, multimedia	<b>Data in Doubt</b> Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations

# Harnessing Big Data



- **OLTP**: Online Transaction Processing (DBMSs)
- **OLAP**: Online Analytical Processing (Data Warehousing)
- **RTAP**: Real-Time Analytics Processing (Big Data Architecture & technology)

# The Model Has Changed...

## The Model of Generating/Consuming Data has Changed

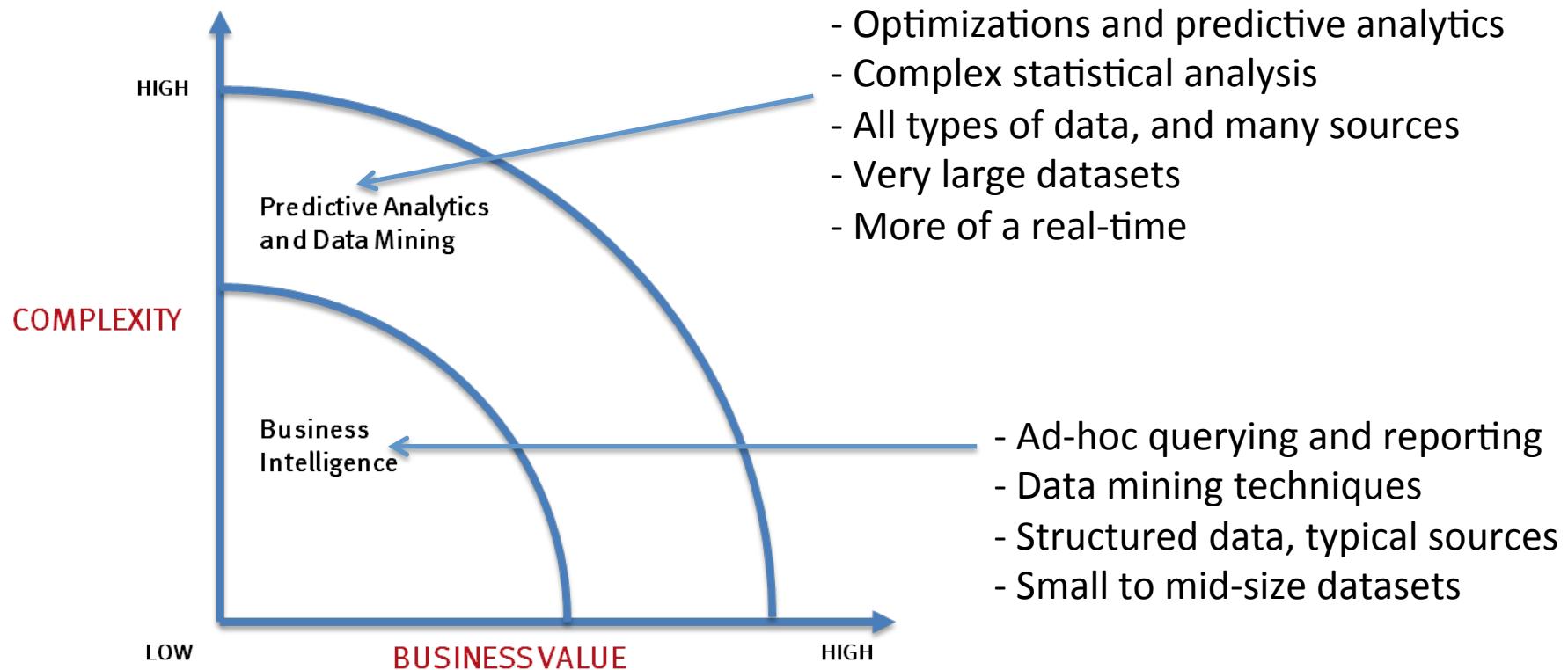
**Old Model:** Few companies are generating data, all others are consuming data



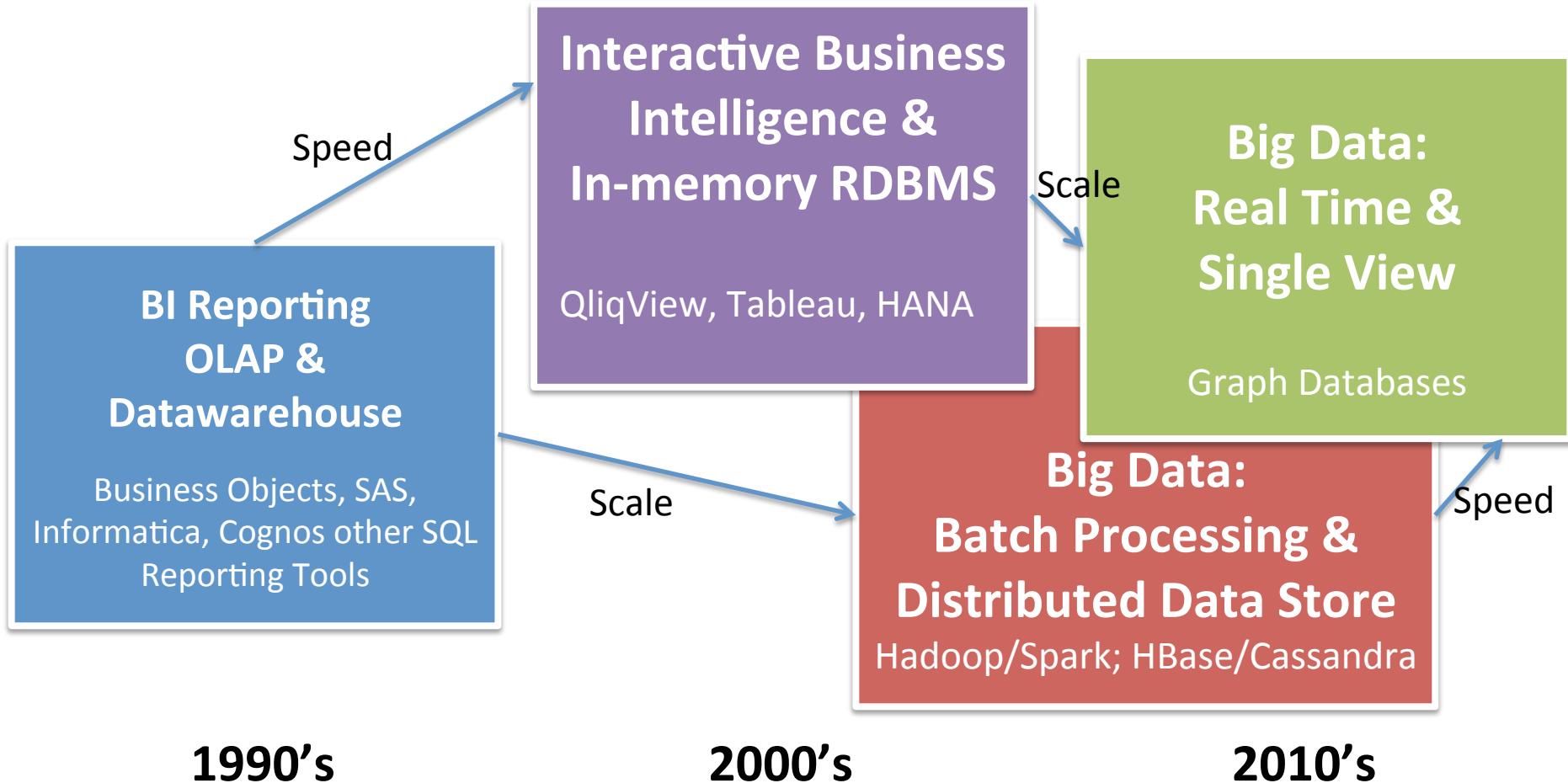
**New Model:** all of us are generating data, and all of us are consuming data



# What's Driving Big Data

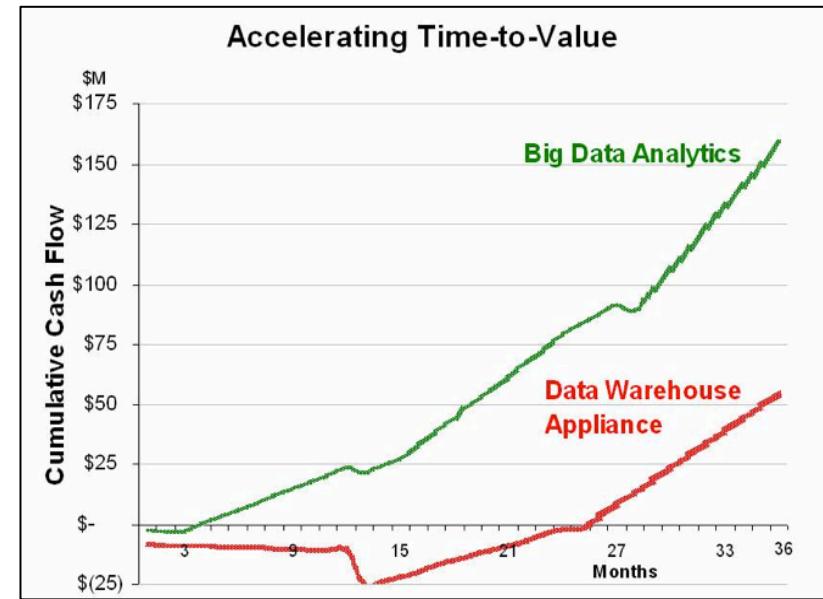


# The Evolution of Business Intelligence



# Big Data Analytics

- Big data is more real-time in nature than traditional DW applications
- Traditional DW architectures (e.g. Exadata, Teradata) are not well-suited for big data apps
- Shared nothing, massively parallel processing, scale out architectures are well-suited for big data apps



# The Big Data Landscape

## Apps

### Vertical Apps



### Operational Intelligence



### Data As A Service



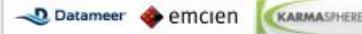
### Ad / Media Apps



### Business Intelligence



### Analytics And Visualization



## Infrastructure

### Analytics Infrastructure



### Operational Infrastructure



### Infrastructure As A Service



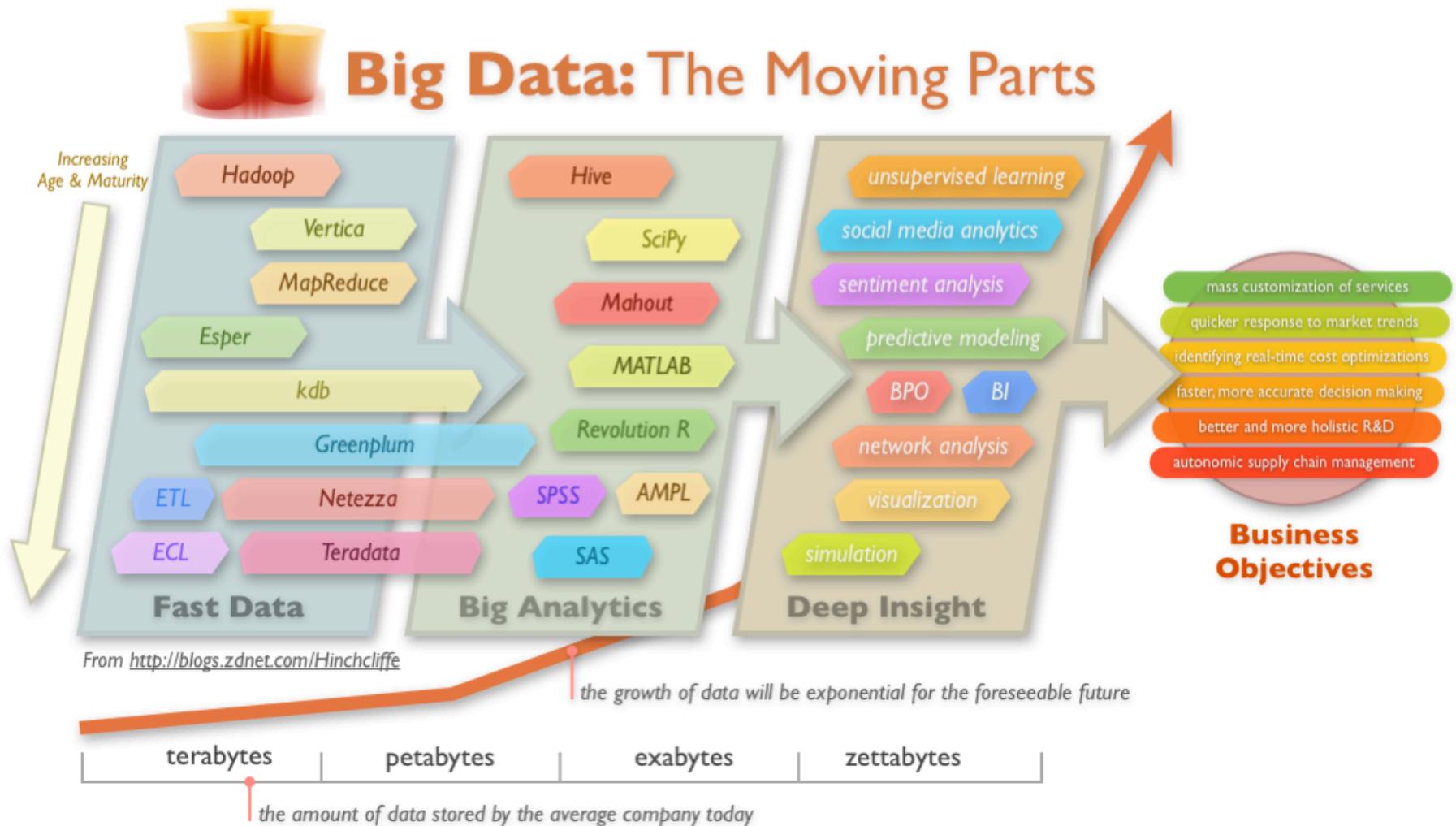
### Structured Databases



## Technologies



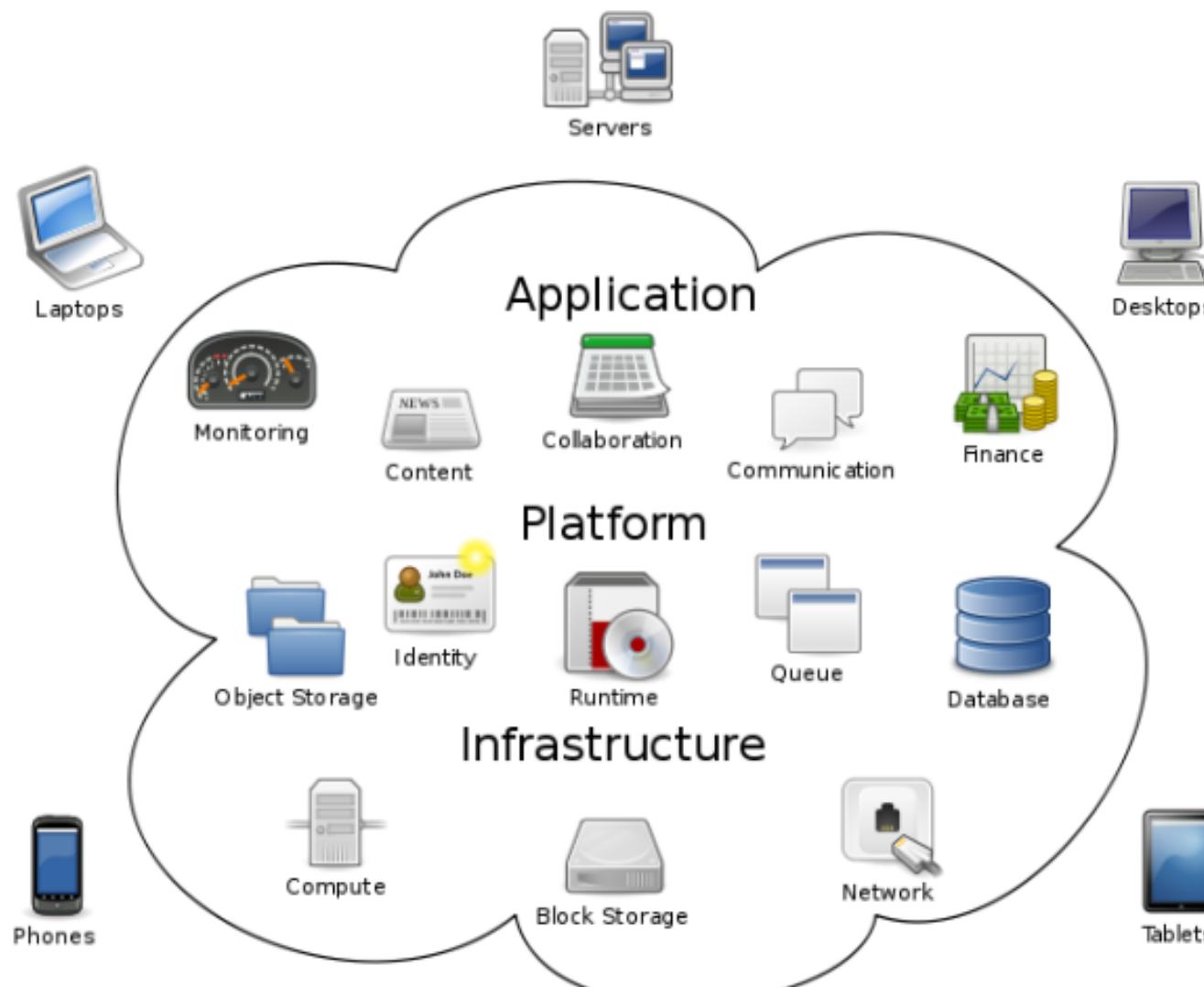
# Big Data Technology



# Cloud Computing

- IT resources provided as a service
  - Compute, storage, databases, queues
- Clouds leverage economies of scale of commodity hardware
  - Cheap storage, high bandwidth networks & multicore processors
  - Geographically distributed data centers
- Offerings from Microsoft, Amazon, Google, ...

# Cloud Computing

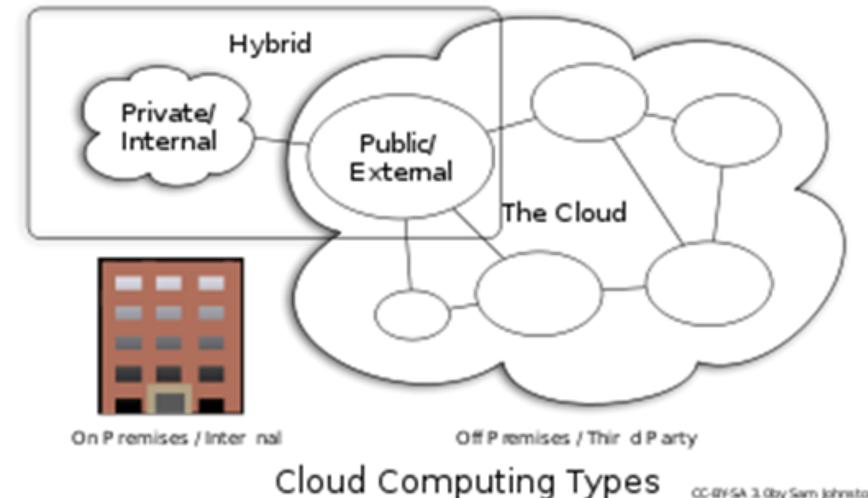


# Benefits

- Cost & management
  - Economies of scale, “out-sourced” resource management
- Reduced Time to deployment
  - Ease of assembly, works “out of the box”
- Scaling
  - On demand provisioning, co-locate data and compute
- Reliability
  - Massive, redundant, shared resources
- Sustainability
  - Hardware not owned

# Types of Cloud Computing

- **Public Cloud:** Computing infrastructure is hosted at the vendor's premises.
- **Private Cloud:** Computing architecture is dedicated to the customer and is not shared with other organisations.
- **Hybrid Cloud:** Organisations host some critical, secure applications in private clouds. The not so critical applications are hosted in the public cloud
- **Cloud Bursting:** The organisation uses its own infrastructure for normal usage, but cloud is used for peak loads.
- **Community Cloud**

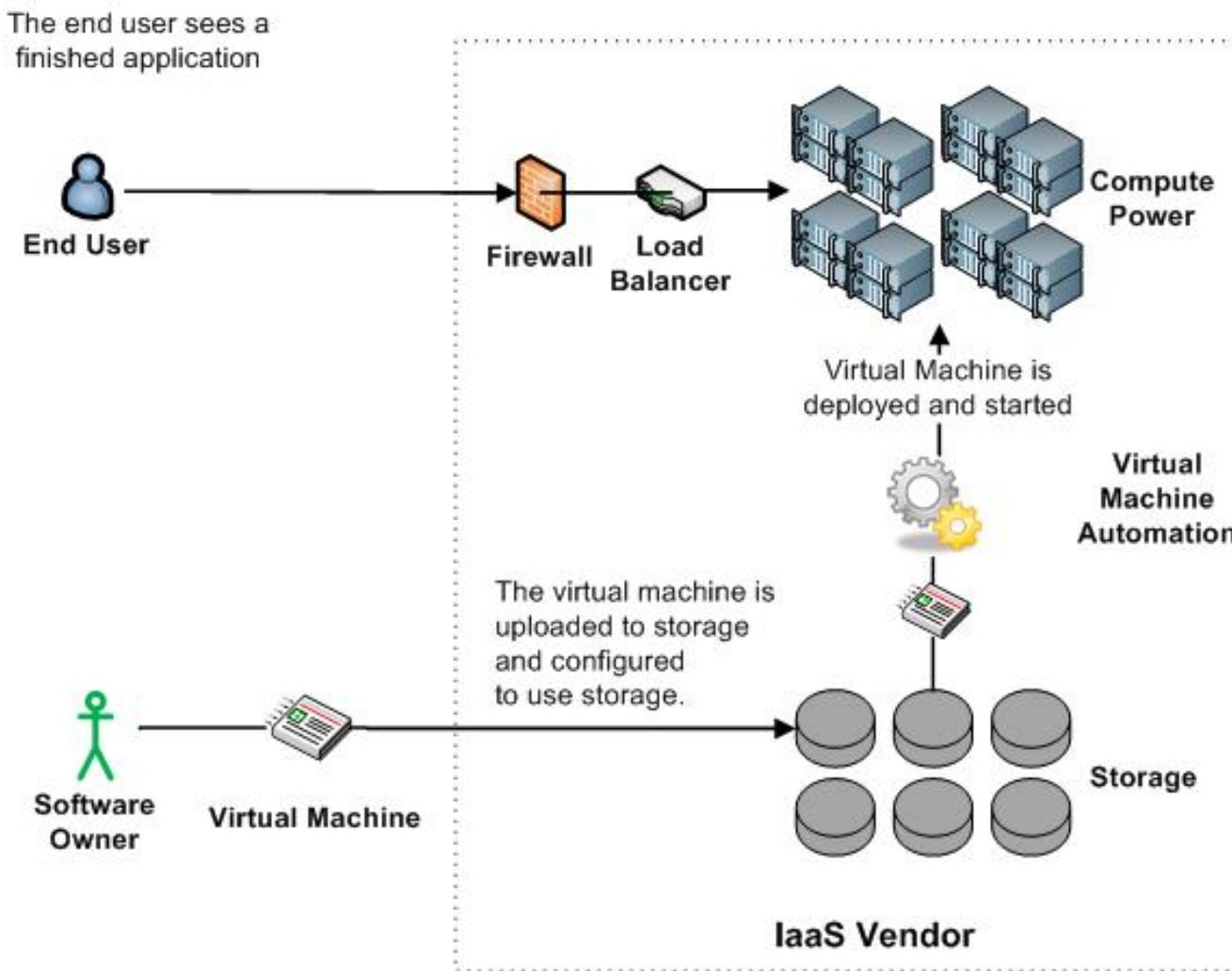


CC-BY-SA 3.0 by Sam Johnston

# Classification of Cloud Computing Based on Service Provided

- Infrastructure as a service (IaaS)
  - Offering hardware related services using the principles of cloud computing. These could include storage services (database or disk storage) or virtual servers.
  - [Amazon EC2](#), [Amazon S3](#), [Rackspace Cloud Servers](#) and [Flexiscale](#).
- Platform as a Service (PaaS)
  - Offering a development platform on the cloud.
  - [Google's Application Engine](#), [Microsoft's Azure](#), [Salesforce.com's force.com](#).
- Software as a service (SaaS)
  - Including a complete software offering on the cloud. Users can access a software application hosted by the cloud vendor on pay-per-use basis. This is a well-established sector.
  - Salesforce.com's offering in the online Customer Relationship Management (CRM) space, Google's [Gmail](#) and Microsoft's [Hotmail](#), [Google Docs](#).

# Infrastructure as a Service (IaaS)



# More Refined Categorization

- Storage-as-a-service
- Database-as-a-service
- Information-as-a-service
- Process-as-a-service
- Application-as-a-service
- Platform-as-a-service
- Integration-as-a-service
- Security-as-a-service
- Management/  
Governance-as-a-service
- Testing-as-a-service
- Infrastructure-as-a-service

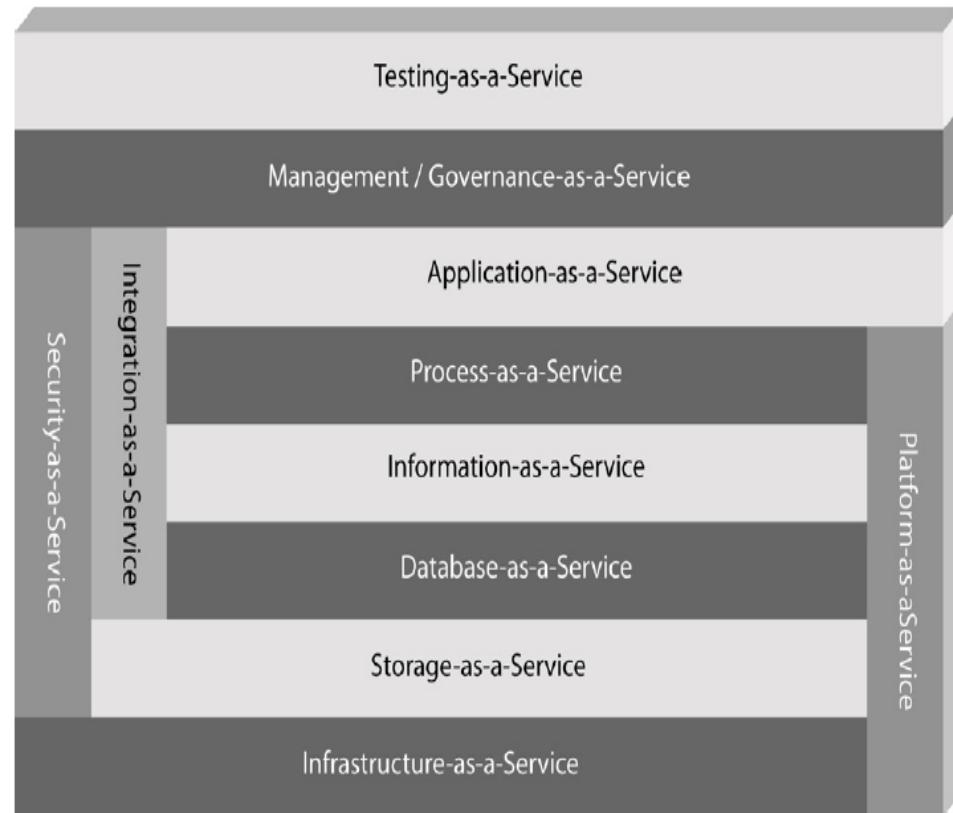


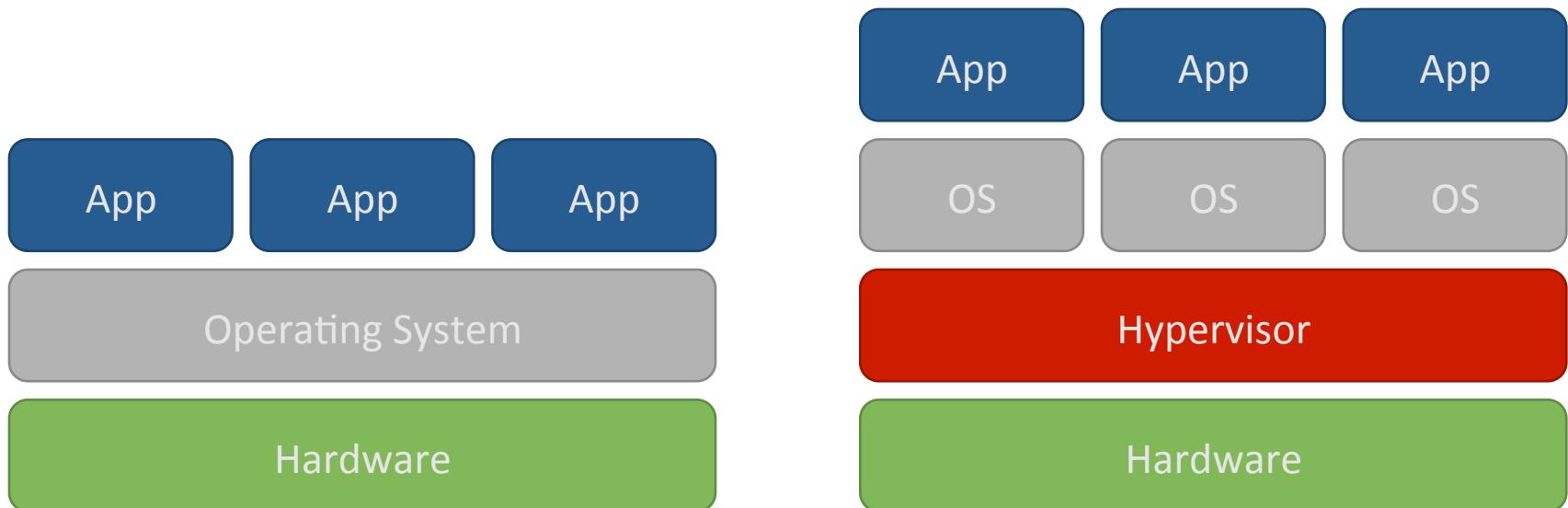
Figure 1: The patterns or categories of cloud computing providers allow you to use a discrete set of services within your architecture.

InfoWorld Cloud Computing Deep Dive

# Key Ingredients in Cloud Computing

- Service-Oriented Architecture (SOA)
- Utility Computing (on demand)
- Virtualization
- SAAS (Software As A Service)
- PAAS (Platform As A Service)
- IAAS (Infrastructure As A Service)
- Web Services in Cloud

# Enabling Technology: Virtualization



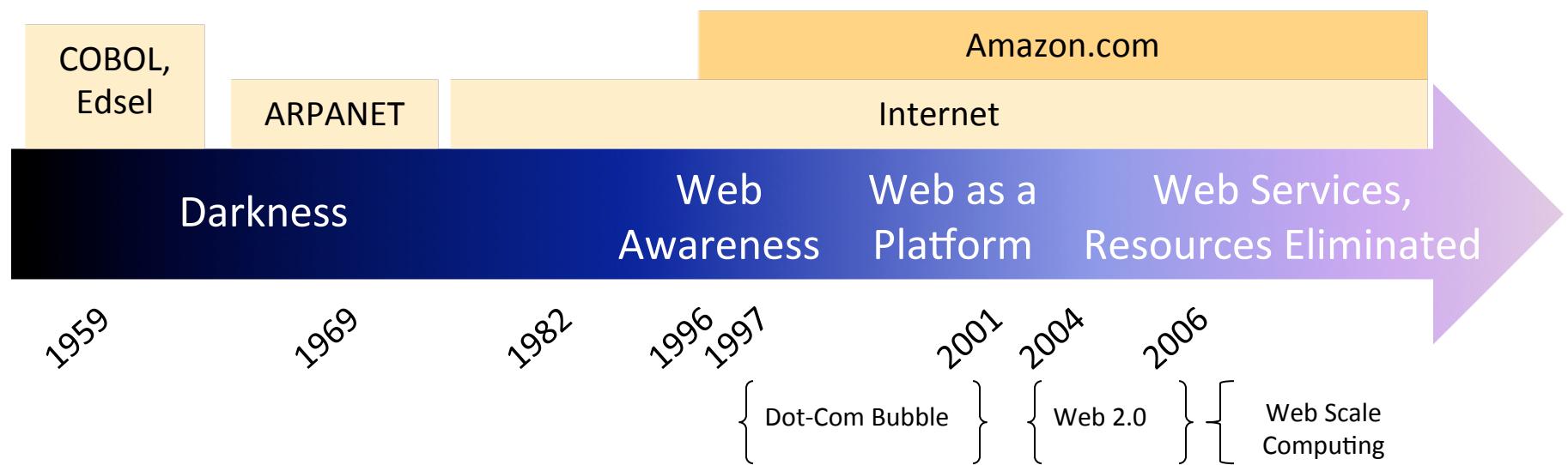
# Everything as a Service

- Utility computing = Infrastructure as a Service (IaaS)
  - Why buy machines when you can rent cycles?
  - Examples: Amazon's EC2, Rackspace
- Platform as a Service (PaaS)
  - Give me nice API and take care of the maintenance, upgrades, ...
  - Example: Google App Engine
- Software as a Service (SaaS)
  - Just run it for me!
  - Example: Gmail, Salesforce

# Cloud versus Cloud

- Amazon Elastic Compute Cloud
- Google App Engine
- Microsoft Azure
- GoGrid
- AppNexus

# The Obligatory Timeline Slide



# Amazon Web Services

- Elastic Compute Cloud – EC2 (IaaS)
- Simple Storage Service – S3 (IaaS)
- Elastic Block Storage – EBS (IaaS)
- SimpleDB (SDB) (PaaS)
- Simple Queue Service – SQS (PaaS)
- CloudFront (S3 based Content Delivery Network – PaaS)
- Consistent AWS Web Services API

# What does Azure platform offer to developers?

## Your Applications



Service Bus

Workflow

Access Control

...



Database

Analytics

Reporting

...



Live Services

Identity

Contacts

Devices

...

...

Compute

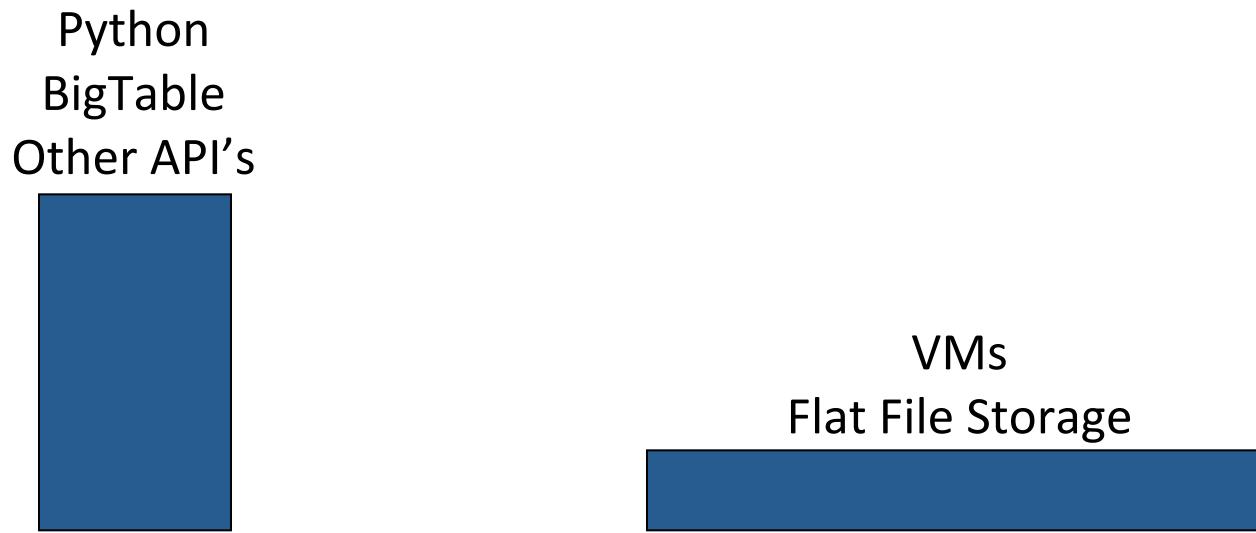
Storage

Manage

...



# Google's AppEngine vs. Amazon's EC2



## AppEngine:

- Higher-level functionality (e.g., automatic scaling)
- More restrictive (e.g., respond to URL only)
- Proprietary lock-in

## EC2/S3:

- Lower-level functionality
- More flexible
- Coarser billing model

# Textbooks

- No Official Textbooks
- References:
  - Hadoop: The Definitive Guide, Tom White, O'Reilly
  - Hadoop In Action, Chuck Lam, Manning
  - Data-Intensive Text Processing with MapReduce, Jimmy Lin and Chris Dyer (  
[www.umiacs.umd.edu/~jimmylin/MapReduce-book-final.pdf](http://www.umiacs.umd.edu/~jimmylin/MapReduce-book-final.pdf))
  - Data Mining: Concepts and Techniques, Third Edition, by Jiawei Han et al.
- Many Online Tutorials and Papers

# Cloud Resources

- Hadoop on your local machine
- Hadoop in a virtual machine on your local machine (**Pseudo-Distributed on Ubuntu**)
- Hadoop in the clouds with Amazon EC2
- Please sign up for:

<https://aws.amazon.com/education/awseducate/>